

**CLAIMS**

1 1. In a cluster of computing nodes having shared access  
2 to one or more volumes of data storage using a parallel  
3 file system, a method for managing the data storage,  
4 comprising:

5 initiating a session of a data management  
6 application on a session node selected from among the  
7 nodes in the cluster;

8 receiving an event message in a session queue for  
9 processing by the data management application at the  
10 session node, responsive to a request submitted to the  
11 parallel file system by a user application on a source  
12 node among the nodes in the cluster to perform a file  
13 operation on a file in the data storage; and

14 following a failure at the session node,  
15 reconstructing the session queue so that processing of  
16 the event message by the data management application can  
17 continue after recovery from the failure.

1 2. A method according to claim 1, wherein the failure  
2 at the session node comprises a file system failure at  
3 the session node.

1 3. A method according to claim 2, wherein  
2 reconstructing the session queue comprises selecting a  
3 new session node from among the nodes on which the file  
4 system failure has not occurred, and moving the data  
5 management session to the new session node.

1 4. A method according to claim 1, wherein  
2 reconstructing the session queue comprises selecting a  
3 new session node from among the nodes in the cluster, and  
4 assuming the data management session on the new session

5 node, whereupon the session queue is reconstructed on the  
6 new session node.

1 5. A method according to claim 4, wherein assuming the  
2 data management session comprises moving the session to a  
3 different node from the session node used before the  
4 failure.

1 6. A method according to claim 4, wherein assuming the  
2 data management session comprises assuming the session on  
3 the same session node that was used before the failure.

1 7. A method according to claim 6, wherein assuming the  
2 session comprises explicitly invoking a session creation  
3 function call of a data management application  
4 programming interface (DMAPI).

1 8. A method according to claim 6, wherein the failure  
2 comprises a file system failure at the session node,  
3 which is followed by file system recovery, and wherein  
4 assuming the session comprises invoking any function call  
5 of a data management application programming interface  
6 (DMAPI) at the session node after the recovery, whereby  
7 reconstruction of the session queue is triggered  
8 implicitly.

1 9. A method according to claim 1, and comprising  
2 storing information regarding the session and events  
3 before the failure at one or more additional nodes among  
4 the nodes in the cluster, wherein reconstructing the  
5 session queue comprises using the information stored at  
6 the one or more additional nodes to reconstruct the  
7 queue.

1 10. A method according to claim 1, and comprising  
2 selecting one of the nodes to serve as a session manager

3 node, and assuming the session by sending a message to  
4 the session manager node, causing the session manager  
5 node to distribute information regarding the session  
6 among the nodes in the cluster so that the data  
7 management application can continue after the recovery.

1 11. A method according to claim 1, wherein initiating  
2 the session comprises initiating the session in  
3 accordance with a data management application programming  
4 interface (DMAPI) of the parallel file system, and  
5 wherein processing the event message comprises processing  
6 the request using the DAPI.

1 12. A method according to claim 1, and comprising:  
2 sending a response to the event message from the  
3 data management application on the session node to the  
4 source node following the recovery from the failure; and  
5 performing the file operation requested by the  
6 source node subject to the response from the data  
7 management application.

1 13. A method according to claim 12, wherein receiving  
2 the event message comprises receiving the message  
3 responsive to submission of the request by a file  
4 operation thread of a user application running on the  
5 source node, and blocking the thread until the response  
6 is received from the session node after the recovery from  
7 the failure.

1 14. A method according to claim 13, wherein  
2 reconstructing the session queue comprises sending a  
3 message from the session node to all of the nodes, so as  
4 to prompt the file operation thread on the source node to  
5 submit a new event message to the session node, whereby

6 the event is placed in the reconstructed queue responsive  
7 to the new message.

1 15. A method according to claim 14, wherein prompting  
2 the file operation thread comprises instructing the file  
3 operation thread to submit the new event message with  
4 respect to an event that is defined as a synchronous  
5 event.

1 16. A method according to claim 15, wherein an event  
2 that is defined as an asynchronous event that was in the  
3 session queue prior to the failure is not placed in the  
4 reconstructed queue.

1 17. A method according to claim 14, wherein receiving  
2 the event message comprises receiving an event  
3 identifier, which is assigned to the event at the source  
4 node, and wherein the event placed in the reconstructed  
5 queue has the same event identifier as was assigned  
6 before the failure.

1 18. A method according to claim 1, and comprising  
2 processing the event message in the reconstructed queue,  
3 and responsive to the event message, reacquiring a data  
4 management access right needed to handle the request.

1 19. A method according to claim 1, wherein receiving the  
2 event message comprises receiving multiple event messages  
3 from multiple source nodes in the cluster, and wherein  
4 reconstructing the session queue comprises collecting  
5 information regarding the session and events from the  
6 multiple source nodes.

1 20. A method according to claim 1, wherein initiating  
2 the session of the data management application comprises  
3 initiating a data migration application, so as to free

4 storage space on at least one of the volumes of data  
5 storage.

1 21. A method according to claim 1, and comprising,  
2 following the failure, when the source node has not  
3 received a response to the event message within a  
4 predetermined lapse of time, failing the request  
5 submitted at the source node to the parallel file system.

1 22. Computing apparatus, comprising:

2 one or more volumes of data storage, arranged to  
3 store data; and

4 a plurality of computing nodes, linked to access the  
5 volumes of data storage using a parallel file system, and  
6 arranged so as to enable a data management application to  
7 initiate a data management session on a session node  
8 selected among the nodes in the cluster, so that when a  
9 request is submitted to the parallel file system by a  
10 user application on a source node among the nodes in the  
11 cluster to perform a file operation on a file in the data  
12 storage, an event message is received at the session node  
13 responsive to the request, for processing by the data  
14 management application, and so that following a failure  
15 at the session node, the session queue is reconstructed  
16 so that processing of the event message by the data  
17 management application can continue after recovery from  
18 the failure.

1 23. Apparatus according to claim 22, wherein the failure  
2 at the session node comprises a file system failure at  
3 the session node.

1 24. Apparatus according to claim 23, wherein the nodes  
2 are arranged so that following the file system failure, a

3 new session node is selected from among the nodes on  
4 which the file system failure has not occurred, and the  
5 data management session is moved to the new session node.

1 25. Apparatus according to claim 22, wherein the nodes  
2 are arranged so that following the failure, a new session  
3 node is selected from among the nodes on which the  
4 failure has not occurred, and the data management session  
5 is assumed on the new session node, whereupon the session  
6 queue is reconstructed on the new session node.

1 26. Apparatus according to claim 25, wherein the session  
2 is assumed on a different node from the session node used  
3 before the failure.

1 27. Apparatus according to claim 25, wherein the session  
2 is assumed on the same session node that was used before  
3 the failure.

1 28. Apparatus according to claim 27, wherein the session  
2 is assumed by explicitly invoking a session creation  
3 function call of a data management application  
4 programming interface (DMAPI).

1 29. Apparatus according to claim 27, wherein the failure  
2 comprises a file system failure at the session node,  
3 which is followed by file system recovery, and wherein  
4 the session is assumed by invoking any function call of a  
5 data management application programming interface (DMAPI)  
6 at the session node after the recovery, whereby  
7 reconstruction of the session queue is triggered  
8 implicitly.

1 30. Apparatus according to claim 22, wherein the nodes  
2 are arranged so that information regarding the session  
3 and events is stored before the failure at one or more

4 additional nodes among the nodes in the cluster, whereby  
5 the session queue is reconstructed using the information  
6 stored at the one or more additional nodes.

1 31. Apparatus according to claim 22, wherein one of the  
2 nodes is selected to serve as a session manager node, and  
3 wherein to assume the session, a message is sent to the  
4 session manager node, causing the session manager node to  
5 distribute information regarding the session among the  
6 nodes in the cluster so that the data management  
7 application can continue after the recovery.

1 32. Apparatus according to claim 22, wherein the session  
2 is initiated in accordance with a data management  
3 application programming interface (DMAPI) of the parallel  
4 file system, and wherein the event message is processed  
5 using the DMAPI.

1 33. Apparatus according to claim 22, wherein the nodes  
2 are arranged so that a response to the event message is  
3 sent from the data management application on the session  
4 node to the source node following the recovery from the  
5 failure, whereupon the file operation requested by the  
6 source node is carried out subject to the response from  
7 the data management application.

1 34. Apparatus according to claim 33, wherein the event  
2 message is received responsive to submission of the  
3 request by a file operation thread of a user application  
4 running on the source node, and the thread is blocked  
5 until the response is received from the session node  
6 after the recovery from the failure.

1 35. Apparatus according to claim 34, wherein to  
2 reconstruct the session queue, a message is sent from the

3 session node to all of the nodes, so that the file  
4 operation thread on the source node is prompted to submit  
5 a new event message to the session node, whereby the  
6 event is placed in the reconstructed queue responsive to  
7 the new message.

1 36. Apparatus according to claim 35, wherein the file  
2 operation thread is prompted to submit the new event  
3 message with respect to an event that is defined as a  
4 synchronous event.

1 37. Apparatus according to claim 36, wherein an event  
2 that is defined as an asynchronous event that was in the  
3 session queue prior to the failure is not placed in the  
4 reconstructed queue.

1 38. Apparatus according to claim 35, wherein the event  
2 message contains an event identifier, which is assigned  
3 to the event at the source node, and wherein the event  
4 placed in the reconstructed queue has the same event  
5 identifier as was assigned before the failure.

1 39. Apparatus according to claim 22, wherein after  
2 reconstructing the session queue, the data management  
3 application reacquires a data management access right  
4 needed to handle the request.

1 40. Apparatus according to claim 22, wherein the nodes  
2 are arranged so that the session node receives multiple  
3 event messages from multiple source nodes in the cluster,  
4 and so that in reconstructing the session queue,  
5 information regarding the session and events is collected  
6 from the multiple source nodes.

1 41. Apparatus according to claim 22, wherein the data  
2 management application comprises a data migration

3 application, for freeing storage space on at least one of  
4 the volumes of data storage.

1 42. Apparatus according to claim 22, wherein following  
2 the failure, when the source node has not received a  
3 response to the event message within a predetermined  
4 lapse of time, the request submitted at the source node  
5 to the parallel file system is failed.

1 43. A computer software product for use in a cluster of  
2 computing nodes having shared access to one or more  
3 volumes of data storage using a parallel file system, the  
4 product comprising a computer-readable medium in which  
5 program instructions are stored, which instructions, when  
6 read by the computing nodes, cause a session of a data  
7 management application to be initiated on a session node  
8 selected among the nodes in the cluster, such that when a  
9 user application on a source node among the nodes in the  
10 cluster submits a request to the parallel file system to  
11 perform a file operation on a file in the data storage,  
12 an event message is received at the session node, for  
13 processing by the data management application, and such  
14 that following a failure at the session node, the session  
15 queue is reconstructed so that processing of the event  
16 message by the data management application can continue  
17 after recovery from the failure.

1 44. A product according to claim 43, wherein the failure  
2 at the session node comprises a file system failure at  
3 the session node.

1 45. A product according to claim 44, wherein following  
2 the file system failure, the instructions cause a new  
3 session node to be selected from among the nodes on which

4 the file system failure has not occurred, whereupon the  
5 data management session is moved to the new session node.

1 46. A product according to claim 43, wherein following  
2 the failure, the instructions cause a new session node to  
3 be selected from among the nodes on which the failure has  
4 not occurred, whereupon the data management session is  
5 assumed on the new session node, and the session queue is  
6 reconstructed on the new session node.

1 47. A product according to claim 46, wherein the session  
2 is assumed on a different node from the session node used  
3 before the failure.

1 48. A product according to claim 46, wherein the session  
2 is assumed on the same session node that was used before  
3 the failure.

1 49. A product according to claim 48, wherein the session  
2 is assumed by explicitly invoking a session creation  
3 function call of a data management application  
4 programming interface (DMAPI).

1 50. A product according to claim 48, wherein the failure  
2 comprises a file system failure at the session node,  
3 which is followed by file system recovery, and wherein  
4 the session is assumed by invoking any function call of a  
5 data management application programming interface (DMAPI)  
6 at the session node after the recovery, whereby  
7 reconstruction of the session queue is triggered  
8 implicitly.

1 51. A product according to claim 43, wherein the  
2 instructions cause information regarding the session and  
3 events to be stored before the failure at one or more  
4 additional nodes among the nodes in the cluster, whereby

5 the session queue is reconstructed using the information  
6 stored at the one or more additional nodes.

1 52. A product according to claim 43, wherein one of the  
2 nodes is selected to serve as a session manager node, and  
3 wherein to assume the session, a message is sent to the  
4 session manager node, causing the session manager node to  
5 distribute information regarding the session among the  
6 nodes in the cluster so that the data management  
7 application can continue after the recovery.

1 53. A product according to claim 43, wherein the product  
2 comprises a data management application programming  
3 interface (DMAPI) of the parallel file system, and  
4 wherein the event message is processed using the DMAPI.

1 54. A product according to claim 43, wherein the  
2 instructions cause a response to the event message to be  
3 sent from the data management application on the session  
4 node to the source node following the recovery from the  
5 failure, whereupon the file operation requested by the  
6 source node is carried out subject to the response from  
7 the data management application.

1 55. A product according to claim 54, wherein the event  
2 message is received responsive to submission of the  
3 request by a file operation thread of a user application  
4 running on the source node, and the thread is blocked  
5 until the response is received from the session node  
6 after the recovery from the failure.

1 56. A product according to claim 55, wherein to  
2 reconstruct the session queue, a message is sent from the  
3 session node to all of the nodes, so that the file  
4 operation thread on the source node is prompted to submit

5 a new event message to the session node, whereby the  
6 event is placed in the reconstructed queue responsive to  
7 the new message.

1 57. A product according to claim 56, wherein the file  
2 operation thread is prompted to submit the new event  
3 message with respect to an event that is defined as a  
4 synchronous event.

1 58. A product according to claim 57, wherein an event  
2 that is defined as an asynchronous event that was in the  
3 session queue prior to the failure is not placed in the  
4 reconstructed queue.

1 59. A product according to claim 56, wherein the event  
2 message contains an event identifier, which is assigned  
3 to the event at the source node, and wherein the event  
4 placed in the reconstructed queue has the same event  
5 identifier as was assigned before the failure.

1 60. A product according to claim 43, wherein after  
2 reconstructing the session queue, the data management  
3 application reacquires a data management access right  
4 needed to handle the request.

1 61. A product according to claim 43, wherein the  
2 instructions cause the session node to receive multiple  
3 event messages from multiple source nodes in the cluster,  
4 and to reconstruct the session queue by collecting  
5 information regarding the session and events from the  
6 multiple source nodes.

7 62. A product according to claim 43, wherein the data  
8 management application comprises a data migration  
9 application, for freeing storage space on at least one of  
10 the volumes of data storage.

1 63. A product according to claim 43, wherein following  
2 the failure, when the source node has not received a  
3 response to the event message within a predetermined  
4 lapse of time, the request submitted at the source node  
5 to the parallel file system is failed.